

LANGUAGE REPRESENTATIONS CAN BE WHAT RECOMMENDERS NEED: FINDINGS AND POTENTIALS

Leheng Sheng¹ An Zhang^{1*} Yi Zhang² Yuxin Chen¹ Xiang Wang² Tat-Seng Chua¹

¹National University of Singapore ²University of Science and Technology of China
 leheng.sheng@u.nus.edu, anzhang@u.nus.edu, zy1230@mail.ustc.edu.cn,
 e1143404@u.nus.edu, xiangwang1223@gmail.com, dcscts@nus.edu.sg

ABSTRACT

Recent studies empirically indicate that language models (LMs) encode rich world knowledge beyond mere semantics, attracting significant attention across various fields. However, in the recommendation domain, it remains uncertain whether LMs implicitly encode user preference information. Contrary to prevailing understanding that LMs and traditional recommenders learn two distinct representation spaces due to the huge gap in language and behavior modeling objectives, this work re-examines such understanding and explores extracting a recommendation space directly from the language representation space. Surprisingly, our findings demonstrate that item representations, when linearly mapped from advanced LM representations, yield superior recommendation performance. This outcome suggests the possible homomorphism between the advanced language representation space and an effective item representation space for recommendation, implying that collaborative signals may be implicitly encoded within LMs. Motivated by the finding of homomorphism, we explore the possibility of designing advanced collaborative filtering (CF) models purely based on language representations without ID-based embeddings. To be specific, we incorporate several crucial components (*i.e.*, a multilayer perceptron (MLP), graph convolution, and contrastive learning (CL) loss function) to build a simple yet effective model, with the language representations of item textual metadata (*i.e.*, title) as the input. Empirical results show that such a simple model can outperform leading ID-based CF models on multiple datasets, which sheds light on using language representations for better recommendation. Moreover, we systematically analyze this simple model and find several key features for using advanced language representations: a good initialization for item representations, superior zero-shot recommendation abilities in new datasets, and being aware of user intention. Our findings highlight the connection between language modeling and behavior modeling, which can inspire both natural language processing and recommender system communities¹.

1 INTRODUCTION

Language models (LMs) have achieved great success across various domains (Vaswani et al., 2017; Devlin et al., 2019; Dubey et al., 2024; OpenAI, 2023), raising a critical question about the knowledge encoded within the language space. Recent studies empirically find that LMs extend beyond semantic understanding to encode comprehensive world knowledge about various domains, such as game states (Li et al., 2023a), lexical attributes (Vulic et al., 2020), and even concepts of space and time (Gurnee & Tegmark, 2023) through language modeling. However, in the domain of recommendation where the integration of LMs is attracting widespread interest (Fan et al., 2023; Li et al., 2023b; Wu et al., 2023a), it remains unclear whether LMs inherently encode relevant information on user preferences and behaviors in the language space.

Currently, one prevailing understanding holds that general LMs and traditional recommenders (*e.g.*, collaborative filtering models (Koren et al., 2009; He et al., 2021)) encode distinct representation

*An Zhang is the corresponding author.

¹Codes are available at <https://github.com/LehengTHU/AlphaRec>

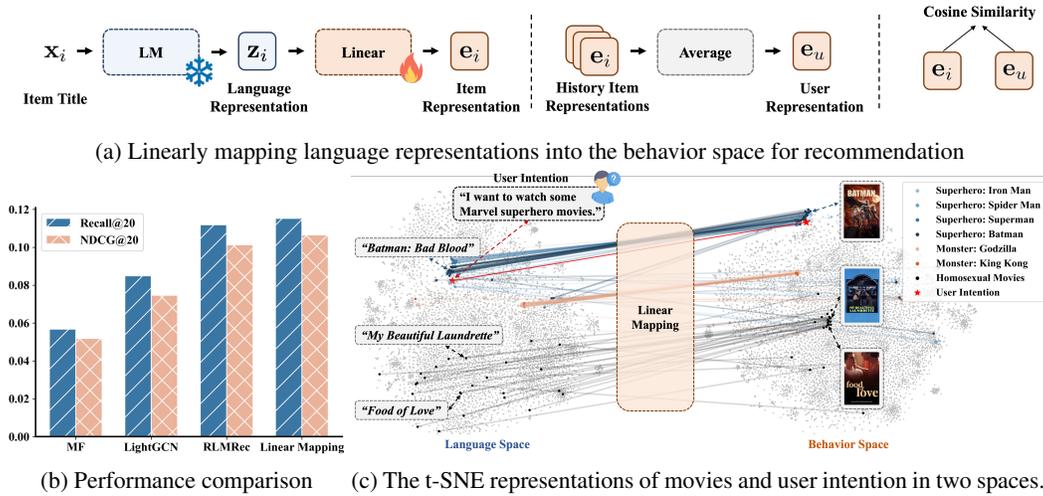


Figure 1: Linearly mapping item titles in language representation space into behavior space yields superior recommendation performance on Movies & TV (Ni et al., 2019) dataset. (1a) The framework of linear mapping. (1b) The recommendation performance comparison between leading CF recommenders and linear mapping. (1c) The t-SNE (Van der Maaten & Hinton, 2008) visualizations of movie representations, with colored lines linking identical movies or user intention across language space (left) and linearly projected behavior space for recommendation (right).

spaces — one for language space and the other for behavior space — but they offer the potential to enhance each other in downstream recommendation tasks (Liao et al., 2024). Specifically, on the one hand, when using LMs as recommenders to directly output items of interest, aligning the language space with the behavior space can significantly improve the recommendation performance (Lin et al., 2023a; Vats et al., 2024; Xu et al., 2024). Various alignment strategies are proposed, including fine-tuning LMs with user behavior data (Zhang et al., 2023d; Bao et al., 2023; Geng et al., 2022; Cui et al., 2022; Lin et al., 2023b), incorporating embeddings from traditional recommenders as a new modality of LMs (Liao et al., 2024; Zhang et al., 2023e; Yang et al., 2023), and extending the vocabulary of LMs with item tokens (Zhu et al., 2023; Zheng et al., 2023; Rajput et al., 2023; Zhai et al., 2024). On the other hand, when using LMs as the enhancer to represent item content (e.g., text metadata), traditional recommenders greatly benefit from text embeddings (Yuan et al., 2022; 2023; Li et al., 2023c; Hou et al., 2024a; Liu et al., 2024a), semantic and reasoning information (Wei et al., 2024; Ren et al., 2024b; Xi et al., 2023), and generated user behaviors (Zhang et al., 2023b;c). Despite these efforts, explicit studies of the relationship between language and behavior spaces remain largely unexplored in the recommendation domain.

In this work, we re-examine this prevailing understanding, by exploring whether LM-generated language space has inherently encoded user preferences and behaviors. Specifically, we test the possibility of directly deriving a behavior space from the language space — that is, we assess whether the language representations of item text metadata (e.g., titles) generated by LMs can independently predict user behaviors and achieve competitive recommendation performance. Positive results would imply that user behavioral patterns, such as collaborative signals (i.e., users’ preference on items reflected by their behavioral similarities) (Wang et al., 2019b), might be implicitly encoded by LMs. To test this hypothesis, we employ linear mapping (Merullo et al., 2023) to project language representations of item titles into a behavior space for recommendation, as Figure 1a shows. Our empirical observations and findings include:

- Before linear mapping, language representation similarities (i.e., semantic textual similarities (STS) (Muennighoff et al., 2023)) may reflect user preference similarities for item contents. Considering Figure 1c as an example, movies with themes of superheroes and monsters cluster together in both language and behavior spaces.
- After linear mapping, language representations are transformed into high-quality behavior representations, which achieve exceptional recommendation performance, as Figure 1b and experimental results in Section 3.2 show. Moreover, the performance improves as the language model size increases and remains relatively robust to prompt disturbances (see Section 3.2).

- Post-mapping language representations encode user behavioral similarities beyond STS. For instance, while certain movies, such as those of homosexual movies (illustrated in Figure 1c), show low STS and their representations disperse in the language space, their projections through linear mapping tend to cluster together, reflecting high user preference similarities.

These findings suggest the **homomorphism** (Dieudonne, 1969) between the LM-generated language space and an expressive behavior space for recommendation. Motivated by this insight, we explore the possibility of building advanced collaborative filtering (CF) models based solely on the language space. To be specific, considering text metadata solely as items’ pre-existing features rather than the widely-used ID information, we perform a frozen LM to create the language representations of items; consequently, we apply a trainable projector (*i.e.*, a two-layer MLP with graph convolution) to map them into a behavior space, and then employ a contrastive loss (*i.e.*, InfoNCE (van den Oord et al., 2018; Wu et al., 2022)) to optimize. We term this model AlphaRec for its simplicity and a series of good properties. Surprisingly, our empirical results show that such a simple model can outperform leading ID-based CF models on multiple datasets. This result sheds light on using language representations for better recommendation.

Furthermore, we systematically analyze this simple model and discover several potentials of adopting language representations for recommendation. First, language representations may serve as a good initialization for item representations, with few adjustments to achieve high recommendation performance (see Section 5.1). This is evidenced by the rapid training convergence of AlphaRec. Second, advanced language representations provide strong zero-shot recommendation capability across entirely new datasets (see Section 5.2). By co-training on multiple datasets, AlphaRec can achieve performance comparable to or even surpassing the fully-trained LightGCN (He et al., 2021) on new datasets without additional training. This underscores the potential of adopting advanced language representations to develop more generalizable recommenders. Third, advanced language representations provide opportunities for perceiving user intentions to refine recommendation results (see Section 5.3). Endowed with the inherent semantic comprehension of language representations, AlphaRec can adjust recommendations according to text-based user intentions, enabling recommenders to evolve into intention-aware systems through a straightforward paradigm shift.

2 PRELIMINARY

2.1 TASK FORMULATION

Personalized recommendation aims to learn user preferences from historical behaviors (*i.e.*, historical interactions with items like view, click, purchase) and find items of interest to trigger users’ future behaviors (Zhou et al., 2018). In this paper, we consider one common recommendation setting: collaborative filtering (CF) (Koren et al., 2022). It aims to select item $i \in \mathcal{I}$ that best matches user u ’s preferences based on binary interaction behaviors $\mathbf{Y} = [y_{ui}]$, where $y_{ui} = 1$ indicates user $u \in \mathcal{U}$ has interacted with item i , and $y_{ui} = 0$ otherwise (Rendle, 2022). Scrutinizing leading CF models, we summarize a common paradigm $\hat{y}_{ui} = s \circ \phi_{\theta}(\mathbf{x}_u, \mathbf{x}_i)$ involving three components:

- For a user-item pair (u, i) , we first get their pre-existing features \mathbf{x}_u and \mathbf{x}_i , which are usually set as ID information or one-hot encodings in CF (Koren et al., 2009; Rendle, 2022; He et al., 2021).
- Upon \mathbf{x}_u and \mathbf{x}_i , the representation generation module ϕ parameterized by θ is adopted to transfer them into behavior representations \mathbf{e}_u and \mathbf{e}_i , encoding the behavioral patterns of users. Its architecture can vary widely, including ID-based embeddings (Koren et al., 2009), multilayer perceptions (He et al., 2017), graph neural networks (Wang et al., 2019b; Cai et al., 2023), and variational autoencoders (Liang et al., 2018).
- Upon \mathbf{e}_u and \mathbf{e}_i , the scoring function s is used to quantify their relevance reflecting how likely user u will interact with item i . One widely-used function is cosine similarity, $s(\mathbf{e}_u, \mathbf{e}_i) = \frac{\mathbf{e}_u^{\top} \mathbf{e}_i}{\|\mathbf{e}_u\| \cdot \|\mathbf{e}_i\|}$ (Chen et al., 2023; Wu et al., 2022).

2.2 ITEM REPRESENTATION GENERATION

Here we emphasize the critical role of item representation generation, which involves transforming item i ’s pre-existing features \mathbf{x}_i into representations \mathbf{e}_i suitable for recommendation. This process is essential, as the quality of these representations directly impacts the recommendation performance.

In this paper, we focus mainly on two kinds of item representation generation tailor-made for different pre-existing features: ID- and LM-based generators.

ID-based generator. Prevailing CF models (Koren et al., 2009; Rendle, 2022; Koren et al., 2022; Wang et al., 2019b; He et al., 2021; Yu et al., 2024) typically convert the ID information of each item i into one-hot encodings (e.g., pre-existing features \mathbf{x}_i). These sparse features are then passed through a trainable generator, such as ID embedding matrices (Koren et al., 2009; Rendle, 2022) or optionally combined with graph convolution layers (Wang et al., 2019b; He et al., 2021; Yu et al., 2024), to generate dense representations \mathbf{e}_i . Optimizing the learning of ID-based representations allows the generator to effectively learn user preferences and behaviors, leading to competitive recommendation performance. However, such ID-based generators suffer several problems, such as poor domain transferability and lack of user intention-aware abilities, since one-hot encodings lack sufficient semantics beyond being identifiers (He et al., 2021).

LM-based generator. Beyond ID information, another research line (Pazzani & Billsus, 2007; Covington et al., 2016; Liu et al., 2024b; Zhang et al., 2024a; Liu et al., 2023b) explores using the text metadata of item i (e.g., titles, descriptions) as pre-existing features \mathbf{x}_i . These features are fed into the LM-based generator, typically a combination of two subsequent components: (1) A frozen LM to extract i 's language representation \mathbf{z}_i , such as the encoder-only LMs like BERT-style models (Devlin et al., 2019; Liu et al., 2019), the decoder-only LMs like Llama-style autoregressive models (Touvron et al., 2023b; Jiang et al., 2023) and OpenAI text embedding models (Neelakantan et al., 2022); and (2) A trainable projector to map \mathbf{z}_i into the final representation \mathbf{e}_i , often using layers like graph convolution layers (He et al., 2021). Although such LM-based generators have been explored to enrich the item representations in literatures (Yuan et al., 2023; Li et al., 2023c; Ren et al., 2024b), few studies have demonstrated that they can solely outperform ID-based generators in recommendation tasks. Worse still, the relationship between the LM-based language space and the behavior space remain largely unexplored in the recommendation domain.

3 UNCOVERING COLLABORATIVE SIGNALS IN LMS VIA LINEAR MAPPING

In this section, we first explore the following research questions. **RQ1:** Do LMs inherently encode collaborative signals (i.e., users' preferences for items as reflected by behavioral similarities) within their representation spaces? **RQ2:** If so, does the presence of such signals scale with model size, and are they robust across different settings? To investigate these questions, we use linear mapping to project language representations of item titles into a behavior space for recommendation. We detail the implementation of the linear mapping in Section 3.1. Subsequently, in Section 3.2, we empirically assess the existence and robustness of collaborative signals in language representations.

3.1 LINEAR MAPPING

Linear mapping is effective to study the representation properties of LMs (Merullo et al., 2023; Alain & Bengio, 2017), discovering the homomorphism (Dieudonne, 1969) between the language space and another space in the target domain. However, its application in the recommendation domain remains largely underexplored.

To bridge this gap, we train a linear mapping matrix \mathbf{W} to project representations from the language space into a behavior space for recommendation. High performance of this linear mapping on the test set would indicate the presence of homomorphism between the language space and an effective behavior space, suggesting the possible existence of collaborative signals in the language representation space (Ravichander et al., 2021; Gurnee & Tegmark, 2023). The overall framework of linear mapping is illustrated in Figure 1a. Specifically, we use frozen LMs to transform pre-existing item title features \mathbf{x}_i into language representations \mathbf{z}_i . To derive user representations, we compute the average of the language representations of the items a user u has interacted with, denoted as $\mathbf{z}_u = \frac{1}{|\mathcal{N}_u|} \sum_{i \in \mathcal{N}_u} \mathbf{z}_i$, where \mathcal{N}_u is the set of user u 's historical items. See Appendix B.2 for detailed procedures for obtaining language representations. The linear mapping matrix sets behavior representations of user u and item i as $\mathbf{e}_u = \mathbf{W}\mathbf{z}_u$ and $\mathbf{e}_i = \mathbf{W}\mathbf{z}_i$ respectively. To optimize the matrix \mathbf{W} , we adopt the InfoNCE loss (van den Oord et al., 2018) as the objective function, which has demonstrated strong performance in both ID-based (Zhang et al., 2023a; Yu et al., 2024) and LM-based generators (Ren et al., 2024b) (refer to equation 4 for the formula).

Table 1: The comparison of the recommendation performance of linear mapping with the classical ID-based CF baselines.

		Movies & TV			Video Games			Books		
		Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
CF	MF (Rendle et al., 2012)	0.0568	0.0519	0.3377	0.0323	0.0195	0.0864	0.0437	0.0391	0.2476
	MultVAE (Liang et al., 2018)	0.0853	0.0776	0.4434	0.0908	0.0531	0.2211	0.0722	0.0597	0.3418
	LightGCN (He et al., 2021)	0.0849	0.0747	0.4397	0.1007	0.0590	0.2281	0.0723	0.0608	0.3489
Linear Mapping	BERT	0.0415	0.0399	0.2362	0.0524	0.0309	0.1245	0.0226	0.0194	0.1240
	RoBERTa	0.0406	0.0387	0.2277	0.0578	0.0338	0.1339	0.0247	0.0209	0.1262
	Llama2-7B	0.1027	0.0955	0.4952	0.1249	0.0729	0.2746	0.0662	0.0559	0.3176
	Mistral-7B	0.1039	0.0963	0.4994	0.1270	0.0687	0.2428	0.0650	0.0544	0.3124
	text-embedding-ada-v2	0.0926	0.0874	0.4563	0.1176	0.0683	0.2579	0.0515	0.0436	0.2570
	text-embeddings-3-large	0.1109	0.1023	0.5200	0.1367	0.0793	0.2928	0.0735	0.0608	0.3355
	SFR-Embedding-Mistral	0.1152	0.1065	0.5327	0.1370	0.0787	0.2927	0.0738	0.0610	0.3371

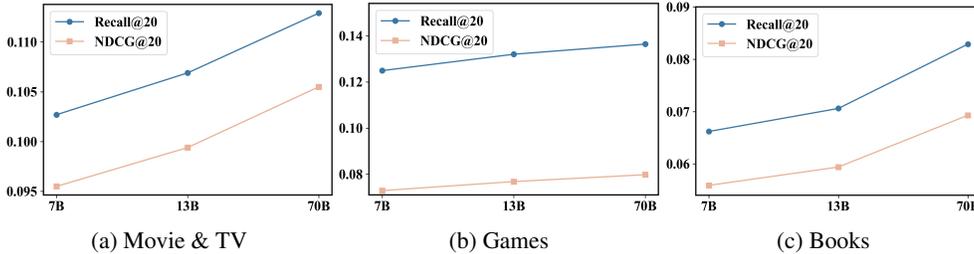


Figure 2: The recommendation performance of linear mapping with different language model sizes.

3.2 EMPIRICAL FINDINGS

Existence (RQ1). To explore the existence of collaborative signals in language representations, we test the recommendation performance of the linear mapping method. Table 1 reports the performance yielded by post-mapping representations on three Amazon datasets (Ni et al., 2019), comparing with classic ID-based CF baselines: matrix factorization (MF) (Koren et al., 2009), MultVAE (Liang et al., 2018), and LightGCN (He et al., 2021) (see baseline details in Appendix C.2). Figures (2a) - (2c) depict the linear mapping performance under different LM sizes. Figure 1c demonstrates the visualization of representations before and after linear mapping. We observe that:

- **Post-mapping representations of advanced LMs achieve superior recommendation performance in most cases, suggesting the possible homomorphism between language spaces and behavior spaces.** Specifically, advanced LMs (e.g., Llama2-7B (Touvron et al., 2023b) and text-embeddings-3-large (Neelakantan et al., 2022)) consistently perform better than leading CF models (e.g., LightGCN) on most metrics. We also see that the performance improves with more recent and advanced LMs. In contrast, earlier BERT-style models (e.g., BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019)) perform similarly to or worse than MF, indicating that LMs have only recently developed the ability to encode user preference similarities effectively.
- **Language representations encode user preference similarities beyond semantic textual similarities (STS).** Consider Figure 1c as an example again, homosexual movies, which differ significantly in textual meaning, cluster together after linear mapping. This suggests that user preferences, which are not immediately apparent from text alone, are implicitly encoded in the language space and can be uncovered through linear mapping.

Scaling and Robustness (RQ2). We also investigate whether the encoding of user preference similarities in LMs scales with model size and whether this encoding is robust in the presence of noise in the input prompts. To this end, we test the linear mapping performance of LMs of various sizes and prompts with noise:

- **The encoding of user preference similarities becomes more refined as model size increases, leading to better linear mapping performance.** Specifically, we test the linear mapping performance across different language model sizes (7B, 13B, and 70B) of the Llama2 family (Touvron et al., 2023b). Llama3 (Dubey et al., 2024) is not selected due to the lack of a 13B model. As shown in Figure 2, linear mapping performance improves consistently as the model size increases from 7B to 70B, indicating that larger models capture more nuanced user behavioral patterns.

Table 2: The robustness of language representations for recommendation.

	Movies & TV			Video Games			Books		
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
Title + Random Noise	0.0952	0.0887	0.4731	0.1213	0.0706	0.2722	0.0632	0.0525	0.3099
Title Only	0.1027	0.0955	0.4952	0.1249	0.0729	0.2746	0.0662	0.0559	0.3176

- **Language representations are relatively robust to prompt disturbances.** Following previous works (Gurnee & Tegmark, 2023), we compare two prompting strategies: using item titles alone (e.g., Castlevania), and adding 5-10 random letters to the titles (e.g., Castlevania sdfhsk). Table 2 shows that adding random noise to the item titles had minimal impact on the linear mapping performance. The prompt noise has more impact on Movies & TV since item titles in this dataset are relatively shorter than others (see Appendix C.1). This finding suggests the relative robustness of the recommendation knowledge encoded in the language representation space.

4 LEVERAGING LANGUAGE REPRESENTATIONS FOR BETTER RECOMMENDATION

This finding of possible space homomorphism (Dieudonne, 1969) and encoded collaborative signals arouse interest in the following questions. **RQ3:** How powerful are such language representations for building advanced CF models that can outperform prevailing ID-based CF methods? To address these questions, in Section 4.1, we aim to develop a simple yet effective CF model termed AlphaRec, which is solely based on language representations and merely incorporates three crucial components in modern CF models. After that, we evaluate its performance in Section 4.2 to demonstrate the capability of advanced language representations for recommendation.

4.1 ALPHAREC

We briefly present how this simple model AlphaRec is designed and trained. It is important to highlight that we center on exploring the power of language representations for recommendation, rather than deliberately inventing new CF mechanisms. Generally, the representation generation architecture $\phi_\theta(\cdot, \cdot)$ is simple, which only contains a two-layer MLP and the basic graph convolution operation. The cosine similarity is used as the similarity function $s(\cdot, \cdot)$, and the contrastive loss InfoNCE (van den Oord et al., 2018; Wu et al., 2022) is adopted for optimization. For simplicity, we adopt text-embeddings-3-large (Neelakantan et al., 2022) for language representation generation by default, for its excellent language understanding and representation capabilities.

Nonlinear projection. We substitute the linear matrix delineated in Section 3 with a nonlinear MLP. Nonlinear transformation helps in excavating more comprehensive preference similarities from the language representation space (see discussions about this in Appendix C.4) (He et al., 2017). Taking the averaged language representations of historical items as the user language representation (i.e., $\mathbf{z}_u = \frac{1}{|\mathcal{N}_u|} \sum_{i \in \mathcal{N}_u} \mathbf{z}_i$), the initial nonlinear transformation operation be formulated as:

$$\mathbf{e}_i^{(0)} = \mathbf{W}_2 \text{LeakyReLU}(\mathbf{W}_1 \mathbf{z}_i + \mathbf{b}_1) + \mathbf{b}_2, \quad \mathbf{e}_u^{(0)} = \mathbf{W}_2 \text{LeakyReLU}(\mathbf{W}_1 \mathbf{z}_u + \mathbf{b}_1) + \mathbf{b}_2. \quad (1)$$

Graph convolution. Graph neural networks (GNNs) show superior effectiveness for recommendation (Wang et al., 2019b), owing to the natural user-item graph structure in recommender systems (Wu et al., 2023b). We employ a minimal graph convolution operation (He et al., 2021) to capture more complicated collaborative patterns from high-order connectivity (Wu et al., 2019) as follows:

$$\mathbf{e}_u^{(k+1)} = \sum_{i \in \mathcal{N}_u} \frac{1}{\sqrt{|\mathcal{N}_u|} \sqrt{|\mathcal{N}_i|}} \mathbf{e}_i^{(k)}, \quad \mathbf{e}_i^{(k+1)} = \sum_{u \in \mathcal{N}_i} \frac{1}{\sqrt{|\mathcal{N}_i|} \sqrt{|\mathcal{N}_u|}} \mathbf{e}_u^{(k)}. \quad (2)$$

The information of connected neighbors is aggregated with a symmetric normalization term $\frac{1}{\sqrt{|\mathcal{N}_u|} \sqrt{|\mathcal{N}_i|}}$. Here \mathcal{N}_u (\mathcal{N}_i) denotes the historical item (user) set that user u (item i) has interacted with. The representations $\mathbf{e}_u^{(0)}$ and $\mathbf{e}_i^{(0)}$ projected from the MLP are used as the input of the first layer. After propagating for K layers, the final behavior representation of a user u (item i) is obtained as the average of representations from each layer:

$$\mathbf{e}_u = \frac{1}{K+1} \sum_{k=0}^K \mathbf{e}_u^{(k)}, \quad \mathbf{e}_i = \frac{1}{K+1} \sum_{k=0}^K \mathbf{e}_i^{(k)}. \quad (3)$$

Table 3: The performance comparison with ID-based CF baselines. The improvement achieved by AlphaRec is significant (p -value $\ll 0.05$).

	Movies & TV			Video Games			Books		
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
MF (Rendle et al., 2012)	0.0568	0.0519	0.3377	0.0323	0.0195	0.0864	0.0437	0.0391	0.2476
MultVAE (Liang et al., 2018)	0.0853	0.0776	0.4434	0.0908	0.0531	0.2211	0.0722	0.0597	0.3418
LightGCN (He et al., 2021)	0.0849	0.0747	0.4397	0.1007	0.0590	0.2281	0.0723	0.0608	0.3489
SGL (Wu et al., 2021)	0.0916	0.0838	0.4680	0.1089	0.0634	0.2449	0.0789	0.0657	0.3734
BC Loss (Zhang et al., 2022)	0.1039	0.0943	0.5037	0.1145	0.0668	0.2561	0.0915	0.0779	0.4045
XSimGCL (Yu et al., 2024)	0.1057	0.0984	0.5128	0.1138	0.0662	0.2550	0.0879	0.0745	0.3918
KAR (Xi et al., 2023)	0.1084	0.1001	0.5134	0.1181	0.0693	0.2571	0.0852	0.0734	0.3834
RLMRec (Ren et al., 2024b)	0.1119	0.1013	0.5301	0.1384	0.0809	0.2997	0.0928	0.0774	0.4092
AlphaRec	0.1221*	0.1144*	0.5587*	0.1519*	0.0894*	0.3207*	0.0991*	0.0828*	0.4185*
Imp.% over the best baseline	6.79%	5.34%	2.27%	9.12%	10.75%	5.40%	9.75%	10.51%	7.01%

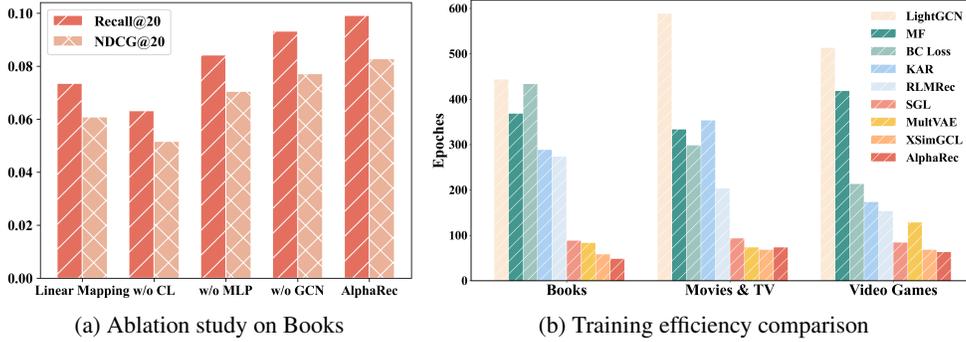


Figure 3: (3a) The effect of each component on Books dataset. (3b) The number of epochs needed for each model to converge. AlphaRec exhibits a breakneck convergence speed.

Contrastive learning objective. The introduction of contrastive learning (Radford et al., 2021) is another key element for the success of leading CF models. Recent research suggests that the contrast learning objective, rather than data augmentation, plays a more significant role in improving recommendation performance (Yu et al., 2024; 2022; Zhang et al., 2023a). Therefore, we simply use the contrast learning object InfoNCE (van den Oord et al., 2018) as the loss function without any additional data augmentation on the graph (Wu et al., 2022). With cosine similarity as the similarity function $s(\mathbf{e}_u, \mathbf{e}_i) = \frac{\mathbf{e}_u^\top \mathbf{e}_i}{\|\mathbf{e}_u\| \cdot \|\mathbf{e}_i\|}$, the InfoNCE loss (van den Oord et al., 2018) is written as:

$$\mathcal{L}_{\text{InfoNCE}} = - \sum_{(u,i) \in \mathcal{O}^+} \log \frac{\exp(s(\mathbf{e}_u, \mathbf{e}_i)/\tau)}{\exp(s(\mathbf{e}_u, \mathbf{e}_i)/\tau) + \sum_{j \in \mathcal{S}_u} \exp(s(\mathbf{e}_u, \mathbf{e}_j)/\tau)}. \quad (4)$$

Here, τ is a hyperparameter called temperature (Wang & Liu, 2021), $\mathcal{O}^+ = \{(u, i) | y_{ui} = 1\}$ denoting the observed interactions between users \mathcal{U} and items \mathcal{I} . And \mathcal{S}_u is a randomly sampled subset of negative items that user u does not adopt.

4.2 EMPIRICAL FINDINGS

Baselines. We compare AlphaRec with leading ID-based CF baselines, to assess the effectiveness of adopting advanced language representations. We do not consider baselines using LMs as recommenders for two practical reasons: the huge inference cost on datasets with millions of interactions and the task limitation of candidate selection (Liao et al., 2024) or next item prediction (Zheng et al., 2023). In addition to classic baselines introduced in section 3.2, we consider two categories of leading ID-based CF baselines, CL-based CF methods: SGL (Wu et al., 2021), BC Loss (Zhang et al., 2022), XSimGCL (Yu et al., 2024) and LM-enhanced CF methods: KAR (Xi et al., 2023), RLMRec (Ren et al., 2024b). See more details about baselines in Appendix C.2.

Recommendation capabilities (RQ3). Table 3 presents performance comparison. The best-performing methods are bold, while the second-best methods are underlined. We observe that:

- **Advanced language representations shows strong potentials for recommendation, which can be unleashed by appropriate model design.** AlphaRec consistently outperforms leading CF

baselines by a large margin across all metrics on all datasets, with an improvement ranging from 6.79% to 9.75% on Recall@20 compared to the best baseline. Moreover, as shown in Figure 3a, each component contributes positively (see more ablation results in Appendix C.3) in unleashing the power of language representations. Specifically, the performance degradation caused by replacing the MLP with a linear weight matrix (w/o MLP) indicates that nonlinear transformations can extract the implicit user preference similarities encoded in the language representation space more effectively. Besides, the performance also drops from replacing InfoNCE loss (Wu et al., 2022) with BPR loss (Rendle et al., 2012) (w/o CL) and removing the graph convolution (w/o GCN) suggests that explicitly modeling the collaborative relationships through the loss function and model architecture can further enhance recommendation performance. These findings suggest that the power of advanced language representations can be unleashed by carefully designing the model, showcasing the potential to surpass prevailing ID-based recommenders.

- **The incorporation of advanced language representations can benefit traditional ID-based CF methods.** We note that two LM-enhanced CF methods, KAR and RLMRec, both show improvements over the most advanced CF methods. Nevertheless, the combination of ID-based embeddings and language representations in these methods does not yield higher results than purely language-representation-based AlphaRec. We attribute this phenomenon to their naive design for the combination ID-based embeddings and language representations, which is also highlighted by previous works (Yuan et al., 2023; Zhang et al., 2024c).

5 EXPLORING POTENTIALS OF LANGUAGE REPRESENTATIONS FOR RECOMMENDATION

In this section, we focus on this question: What new opportunities beyond good performance can advanced language representations bring to recommender systems? To answer this, we systematically analyze AlphaRec and discover the following potentials of adopting such language representations. **Potential 1:** Good initialization for item representations (Section 5.1). **Potential 2:** Zero-shot ability (Section 5.2). **Potential 3:** Intention-aware ability (Section 5.3).

5.1 GOOD INITIALIZATION FOR ITEM REPRESENTATIONS (POTENTIAL 1)

Advanced language representations may provide a good initialization for item representations, with few adjustments for effective recommendation. As shown in Figure 3b, beyond its good performance, AlphaRec also exhibits extremely fast convergence speed, which is comparable with or even surpasses the fastest ID-based CF methods (e.g., SGL (Wu et al., 2021) and XSimGCL (Yu et al., 2024)). Moreover, recent works also suggest that, when using advanced language representations to initialize ID-based item embeddings (Harte et al., 2023; Zhao et al., 2024), the performance of traditional ID-based recommenders improves significantly. We attribute the above findings to the homomorphism between the language space and a good behavior space. Therefore, when using advanced language representations for initialization, only minor adjustments are needed to generate effective behavior representations for recommendation.

5.2 ZERO-SHOT ABILITY (POTENTIAL 2)

The prevailing ID-based recommenders suffer from domain transferring problems (i.e., behavior representations are highly bound with ID information (Zhu et al., 2021)). Advanced language representations may provide opportunities for learning transferable item representations (Hou et al., 2022), enabling recommenders to perform well on entirely new datasets without any ID overlap. To address this potential, we test the zero-shot recommendation ability of AlphaRec (Ding et al., 2021).

Experimental settings. In zero-shot recommendation, there is no item or user overlap between the training set and test set (Ding et al., 2021; Zhang et al., 2024b), which is different from the research line of cross-domain recommendation (Zhu et al., 2021). We jointly train AlphaRec on three source datasets (i.e., Books, Movies & TV, and Video Games), while testing it on three completely new target datasets (i.e., Movielens-1M (Harper & Konstan, 2016), Book Crossing (Lee et al., 2019), and Amazon Industrial & Scientific (Ni et al., 2019)) without further training on these new datasets. (see more details about training on multiple datasets in Appendix D.2.1). Due to the lack of zero-shot recommenders in general CF, we slightly modify the two zero-shot methods in the sequential

Table 4: The zero-shot recommendation performance comparison on entirely new datasets. The improvement achieved by AlphaRec is significant (p -value $\ll 0.05$).

		Industrial & Scientific			MovieLens-1M			Book Crossing		
		Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
full	MF (Rendle et al., 2012)	0.0344	0.0225	0.0521	0.1855	0.3765	0.9634	0.0316	0.0317	0.2382
	MultiVAE (Liang et al., 2018)	0.0751	0.0459	0.1125	0.2039	0.3741	0.9740	0.0736	0.0634	0.3716
	LightGCN (He et al., 2021)	0.0785	0.0533	0.1078	0.2019	0.4017	0.9715	0.0630	0.0588	0.3475
zero-shot	Random	0.0148	0.0061	0.0248	0.0068	0.0185	0.2611	0.0039	0.0036	0.0443
	Pop	0.0216	0.0087	0.0396	0.0253	0.0679	0.5439	0.0119	0.0101	0.1157
	ZESRec (Ding et al., 2021)	0.0326	0.0272	0.0628	0.0274	0.0787	0.5786	0.0155	0.0143	0.1347
	UniSRec (Hou et al., 2022)	0.0453	0.0350	0.0863	0.0578	0.1412	0.7135	0.0396	0.0332	0.2454
	AlphaRec	0.0913*	0.0573	0.1277*	0.1486*	0.3215*	0.9296*	0.0660*	0.0545*	0.3381*
Imp.% over the best zero-shot baseline		157.09%	127.69%	30.29%	66.67%	64.16%	37.78%	101.55%	63.71%	47.97%

recommendation (Wang et al., 2019a), ZESRec (Hou et al., 2022) and UniSRec (Hou et al., 2022), as baselines. We also incorporate two strategy-based CF methods (*i.e.*, Random and Pop) and one method using the large language model (LLM) as zero-shot recommender (*i.e.*, LLMRank (Hou et al., 2024b)) (see more details about baselines in Appendix D.2.2).

Findings. Table 4 presents the zero-shot recommendation performance comparison. The best methods are bold and starred, while the second-best methods are underlined. We observe that:

- **Advanced language representations provide opportunities for learning transferable item representations.** AlphaRec demonstrates strong zero-shot recommendation capabilities, comparable to or even surpassing the fully trained LightGCN. AlphaRec performs better on the Amazon Industrial & Scientific dataset, possibly because it captures user behavioral patterns of the same platform (Ni et al., 2019) through training on multiple Amazon datasets. Conversely, ZESRec and UniSRec exhibit a marked performance decrement compared with AlphaRec. We attribute this phenomenon to two aspects. On the one hand, BERT-style LMs (Devlin et al., 2019; Liu et al., 2019) used in these works may not have effectively encoded user preference similarities, which is consistent with our previous findings in Section 3. On the other hand, components designed for the next item prediction task in sequential recommendation (Kang & McAuley, 2018) may not be suitable for capturing the general preferences of users in CF scenarios. Moreover, AlphaRec also outperforms the leading LLM-based zero-shot recommender LLMRank (see Appendix D.2.3).
- **The zero-shot recommendation capability of advanced language representations generally benefits from an increased amount of training data, without compromising performance on source datasets.** As illustrated in Table 9, the zero-shot performance of AlphaRec, when trained on a mixed dataset, is generally superior to training on one single dataset (Hou et al., 2022). Moreover, we discover that AlphaRec, when trained jointly on multiple datasets, hardly experiences a performance decline on each source dataset. These results indicate the general recommendation capability of a single pre-trained AlphaRec across multiple datasets. The above findings also offer a potential research path to achieve general recommendation capabilities, by incorporating more training data with more themes. See more details about these results in Appendix D.2.4.

5.3 INTENTION-AWARE ABILITY (POTENTIAL 3)

The language understanding ability in advanced language representations (especially representations from LLM-based text embedding models) offers the opportunity for perceiving text-based user intentions and refining recommendations. To study the potential of intention-aware ability, we introduce a new hyperparameter α in AlphaRec to combine user intentions with historical interests.

Experimental settings. To endow AlphaRec with user intention-aware ability, we adopt a simple paradigm shift by introducing a user intention representation $\mathbf{e}_u^{Intention}$. In the inference stage, we obtain the language representation $\mathbf{e}_u^{Intention}$ for each user intention query and combine it with the original user representation to get a new user representation as $\tilde{\mathbf{e}}_u^{(0)} = (1 - \alpha)\mathbf{e}_u^{(0)} + \alpha\mathbf{e}_u^{Intention}$ (Ai et al., 2017). This new user representation $\tilde{\mathbf{e}}_u^{(0)}$ is sent into the pre-trained AlphaRec for recommendation. We test the user intention capture ability of AlphaRec on MovieLens-1M and Video Games. In the test set, only one target item remains for each user (Ai et al., 2017), with one intention query generated by ChatGPT (OpenAI, 2023; Hou et al., 2024a) (see the details about how to generate and check these intention queries in Appendix D.3.1). We report a relatively small $K = 5$ for all metrics to better reflect the intention capture accuracy.

Table 5: The performance comparison in user intention capture.

	MovieLens-1M		Video Games	
	HR@5	NDCG@5	HR@5	NDCG@5
TEM (Bi et al., 2020)	0.2738	0.1973	0.2212	0.1425
AlphaRec (w/o Intention)	0.0793	0.0498	0.0663	0.0438
AlphaRec (w Intention)	0.4704*	0.3738*	0.2569*	0.1862*

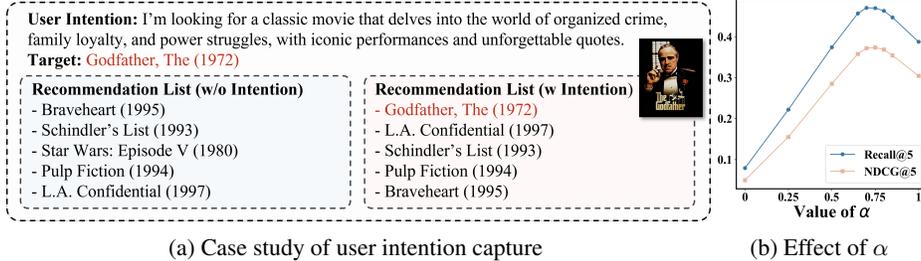


Figure 4: User intention capture experiments on MovieLens-1M. (4a) AlphaRec refines the recommendations according to language-based user intention. (4b) The effect of user intention strength α .

Findings: We report the user intention capture experiment results in Table 5, show one case study in Figure 4a, and study the effect of α in Figure 4b. We find out that:

The language understanding ability in advanced language representations enables recommenders to perceive user intentions and refine recommendations. As shown in Table 5, the introduction of user intention (w Intention) significantly refines the recommendations of the pre-trained AlphaRec (w/o Intention). Moreover, AlphaRec outperforms the baseline model TEM (Bi et al., 2020) by a large margin, even without additional training on search tasks. We further conduct a case study on MovieLens-1M to demonstrate how AlphaRec captures the user intention (see more case study results in Appendix D.3.3). Additionally, the intention-aware ability also benefits from user historical interests. Figure 4b depicts the effect of α . The α controls the strength of user intention, where $\alpha = 0$ denotes that the user intention is neglected and $\alpha = 1$ denotes that the user historical interest is ignored. The convex curve in Figure 4b suggests that both user historical interests and user intention play vital roles. The above findings suggest the potential of adopting advanced language representations to perceive text-based user intentions and refining recommendations (see more details in Appendix D.3).

6 LIMITATIONS

There are several limitations unaddressed in this paper. On the one hand, while we explore the relationship between language spaces and behavior spaces through empirical results, there is no theoretical guarantee. On the other hand, although we investigate the potential of advanced language representations for recommendation, we do not design any new components for CF models.

7 CONCLUSION

In this paper, we explored the relationship between language space and behavior spaces for recommendation, and explored the potential for using language representations for recommendation. Empirical results suggest the possible presence of homomorphism between advanced LMs representation spaces and an effective item representation space for recommendation. Inspired by this finding, we discussed how to unleash the power of advanced language representations by developing a simple yet effective CF model called AlphaRec. Moreover, by systematically analyzing AlphaRec, we explored the potentials of advanced language representations: fast convergence, zero-shot ability, and intention-aware ability. Possible future work will involve exploring the space relationship from both theoretical and multimodal perspectives. We believed that this paper sheds light on re-thinking the connection between language modeling and user behavior modeling, benefiting both natural language processing and recommender system communities.

REFERENCES

- Qingyao Ai, Yongfeng Zhang, Keping Bi, Xu Chen, and W. Bruce Croft. Learning a hierarchical embedding model for personalized product search. In *SIGIR*, 2017.
- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. In *ICLR (Workshop)*, 2017.
- Keqin Bao, Jizhi Zhang, Yang Zhang, Wenjie Wang, Fuli Feng, and Xiangnan He. Tallrec: An effective and efficient tuning framework to align large language model with recommendation. In *RecSys*, 2023.
- Keping Bi, Qingyao Ai, and W. Bruce Croft. A transformer-based embedding model for personalized product search. In *SIGIR*, 2020.
- Xuheng Cai, Chao Huang, Lianghao Xia, and Xubin Ren. Lightgcl: Simple yet effective graph contrastive learning for recommendation. In *ICLR*, 2023.
- Jiawei Chen, Junkang Wu, Jiancan Wu, Xuezhi Cao, Sheng Zhou, and Xiangnan He. Adap- τ : Adaptively modulating embedding magnitude for recommendation. In *WWW*, 2023.
- Yuxin Chen, Junfei Tan, An Zhang, Zhengyi Yang, Leheng Sheng, Enzhi Zhang, Xiang Wang, and Tat-Seng Chua. On softmax direct preference optimization for recommendation. *CoRR*, abs/2406.09215, 2024.
- Zheng Chen. PALR: personalization aware llms for recommendation. *CoRR*, abs/2305.07622, 2023.
- Paul Covington, Jay Adams, and Emre Sargin. Deep neural networks for youtube recommendations. In *RecSys*, 2016.
- Zeyu Cui, Jianxin Ma, Chang Zhou, Jingren Zhou, and Hongxia Yang. M6-rec: Generative pre-trained language models are open-ended recommender systems. *CoRR*, abs/2205.08084, 2022.
- Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu, Zihua Si, Chen Xu, Zhongxiang Sun, Xiao Zhang, and Jun Xu. Uncovering chatgpt’s capabilities in recommender systems. In *RecSys*, 2023.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *ACL*, 2019.
- Jean Dieudonne. *Linear algebra and geometry*. Hermann, 1969.
- Hao Ding, Yifei Ma, Anoop Deoras, Yuyang Wang, and Hao Wang. Zero-shot recommender systems. *CoRR*, abs/2105.08318, 2021.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Rozière, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Grégoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and et al. The llama 3 herd of models. *CoRR*, abs/2407.21783, 2024.

- Wenqi Fan, Zihuai Zhao, Jiatong Li, Yunqing Liu, Xiaowei Mei, Yiqi Wang, Jiliang Tang, and Qing Li. Recommender systems in the era of large language models (llms). *CoRR*, abs/2307.02046, 2023.
- Yunfan Gao, Tao Sheng, Youlin Xiang, Yun Xiong, Haofen Wang, and Jiawei Zhang. Chat-rec: Towards interactive and explainable llms-augmented recommender system. *CoRR*, abs/2303.14524, 2023.
- Binzong Geng, Zhaoxin Huan, Xiaolu Zhang, Yong He, Liang Zhang, Fajie Yuan, Jun Zhou, and Linjian Mo. Breaking the length barrier: Llm-enhanced CTR prediction in long textual user behaviors. *CoRR*, abs/2403.19347, 2024.
- Shijie Geng, Shuchang Liu, Zuohui Fu, Yingqiang Ge, and Yongfeng Zhang. Recommendation as language processing (RLP): A unified pretrain, personalized prompt & predict paradigm (P5). In *RecSys*, 2022.
- Wes Gurnee and Max Tegmark. Language models represent space and time. *CoRR*, abs/2310.02207, 2023.
- F. Maxwell Harper and Joseph A. Konstan. The movielens datasets: History and context. *ACM Trans. Interact. Intell. Syst.*, 5(4):19:1–19:19, 2016.
- Jesse Harte, Wouter Zorgdrager, Panos Louridas, Asterios Katsifodimos, Dietmar Jannach, and Marios Fragkoulis. Leveraging large language models for sequential recommendation. In *RecSys*, 2023.
- Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. Neural collaborative filtering. In *WWW*, 2017.
- Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yong-Dong Zhang, and Meng Wang. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *SIGIR*, 2021.
- Yupeng Hou, Shanlei Mu, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen. Towards universal sequence representation learning for recommender systems. In *KDD*, pp. 585–593. ACM, 2022.
- Yupeng Hou, Zhankui He, Julian J. McAuley, and Wayne Xin Zhao. Learning vector-quantized item representation for transferable sequential recommenders. In Ying Ding, Jie Tang, Juan F. Sequeda, Lora Aroyo, Carlos Castillo, and Geert-Jan Houben (eds.), *WWW*, 2023.
- Yupeng Hou, Jiacheng Li, Zhankui He, An Yan, Xiushi Chen, and Julian J. McAuley. Bridging language and items for retrieval and recommendation. *CoRR*, abs/2403.03952, 2024a.
- Yupeng Hou, Junjie Zhang, Zihan Lin, Hongyu Lu, Ruobing Xie, Julian J. McAuley, and Wayne Xin Zhao. Large language models are zero-shot rankers for recommender systems. In *ECIR*, 2024b.
- Wenyue Hua, Shuyuan Xu, Yingqiang Ge, and Yongfeng Zhang. How to index item ids for recommendation foundation models. In Qingyao Ai, Yiqin Liu, Alistair Moffat, Xuanjing Huang, Tetsuya Sakai, and Justin Zobel (eds.), *SIGIR-AP*, 2023.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L  lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. Mistral 7b. *CoRR*, abs/2310.06825, 2023.
- Wang-Cheng Kang and Julian J. McAuley. Self-attentive sequential recommendation. In *ICDM*, 2018.
- Xiaoyu Kong, Jiancan Wu, An Zhang, Leheng Sheng, Hui Lin, Xiang Wang, and Xiangnan He. Customizing language models with instance-wise lora for sequential recommendation. *CoRR*, abs/2408.10159, 2024.
- Yehuda Koren, Robert M. Bell, and Chris Volinsky. Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37, 2009.

- Yehuda Koren, Steffen Rendle, and Robert M. Bell. Advances in collaborative filtering. In *Recommender Systems Handbook*, pp. 91–142. Springer US, 2022.
- Walid Krichene and Steffen Rendle. On sampled metrics for item recommendation. In *KDD*, 2020.
- Hoyeop Lee, Jinbae Im, Seongwon Jang, Hyunsouk Cho, and Sehee Chung. Melu: Meta-learned user preference estimator for cold-start recommendation. In *KDD*, 2019.
- Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda B. Viégas, Hanspeter Pfister, and Martin Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic task. In *ICLR*, 2023a.
- Lei Li, Yongfeng Zhang, Dugang Liu, and Li Chen. Large language models for generative recommendation: A survey and visionary discussions. *CoRR*, abs/2309.01157, 2023b.
- Ruyu Li, Wenhao Deng, Yu Cheng, Zheng Yuan, Jiaqi Zhang, and Fajie Yuan. Exploring the upper limits of text-based collaborative filtering using large language models: Discoveries and insights. *CoRR*, abs/2305.11700, 2023c.
- Xiangyang Li, Bo Chen, Lu Hou, and Ruiming Tang. Ctrl: Connect tabular and language model for ctr prediction. *arXiv preprint arXiv:2306.02841*, 2023d.
- Dawen Liang, Rahul G. Krishnan, Matthew D. Hoffman, and Tony Jebara. Variational autoencoders for collaborative filtering. In *WWW*, 2018.
- Jiayi Liao, Sihang Li, Zhengyi Yang, Jiancan Wu, Yancheng Yuan, Xiang Wang, and Xiangnan He. Llara: Large language-recommendation assistant. In *SIGIR*, 2024.
- Jianghao Lin, Xinyi Dai, Yunjia Xi, Weiwen Liu, Bo Chen, Xiangyang Li, Chenxu Zhu, Huifeng Guo, Yong Yu, Ruiming Tang, and Weinan Zhang. How can recommender systems benefit from large language models: A survey. *CoRR*, abs/2306.05817, 2023a.
- Jianghao Lin, Rong Shan, Chenxu Zhu, Kounianhua Du, Bo Chen, Shigang Quan, Ruiming Tang, Yong Yu, and Weinan Zhang. Rella: Retrieval-enhanced large language models for lifelong sequential behavior comprehension in recommendation. *CoRR*, abs/2308.11131, 2023b.
- Junling Liu, Chao Liu, Renjie Lv, Kang Zhou, and Yan Zhang. Is chatgpt a good recommender? A preliminary study. *CoRR*, abs/2304.10149, 2023a.
- Qidong Liu, Jiayi Hu, Yutian Xiao, Jingtong Gao, and Xiangyu Zhao. Multimodal recommender systems: A survey. *CoRR*, abs/2302.03883, 2023b.
- Qidong Liu, Xian Wu, Xiangyu Zhao, Yejing Wang, Zijian Zhang, Feng Tian, and Yefeng Zheng. Large language models enhanced sequential recommendation for long-tail user and item. *arXiv preprint arXiv:2405.20646*, 2024a.
- Qijiong Liu, Nuo Chen, Tetsuya Sakai, and Xiao-Ming Wu. ONCE: boosting content-based recommendation with both open- and closed-source large language models. In Luz Angelica Caudillo-Mata, Silvio Lattanzi, Andrés Muñoz Medina, Leman Akoglu, Aristides Gionis, and Sergei Vassilvitskii (eds.), *WSDM*, 2024b.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692, 2019.
- Zhiming Mao, Huimin Wang, Yiming Du, and Kam-Fai Wong. Unitrec: A unified text-to-text transformer and joint contrastive learning framework for text-based recommendation. In *ACL*, 2023.
- Julian J. McAuley, Rahul Pandey, and Jure Leskovec. Inferring networks of substitutable and complementary products. In *KDD*, 2015.

- Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. Sfr-embedding-mistral:enhance text retrieval with transfer learning. Salesforce AI Research Blog, 2024. URL <https://blog.salesforceairesearch.com/sfr-embedded-mistral/>.
- Jack Merullo, Louis Castricato, Carsten Eickhoff, and Ellie Pavlick. Linearly mapping from image to text space. In *ICLR*, 2023.
- Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. MTEB: massive text embedding benchmark. In *EACL*, 2023.
- Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford, Jesse Michael Han, Jerry Tworek, Qiming Yuan, Nikolas Tezak, Jong Wook Kim, Chris Hallacy, Johannes Heidecke, Pranav Shyam, Boris Power, Tyna Eloundou Nekoul, Girish Sastry, Gretchen Krueger, David Schnurr, Felipe Petroski Such, Kenny Hsu, Madeleine Thompson, Tabarak Khan, Toki Sherbakov, Joanne Jang, Peter Welinder, and Lilian Weng. Text and code embeddings by contrastive pre-training. *CoRR*, abs/2201.10005, 2022.
- Jianmo Ni, Jiacheng Li, and Julian J. McAuley. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *EMNLP*, 2019.
- OpenAI. GPT-4 technical report. *CoRR*, 2023.
- Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. *CoRR*, abs/2311.03658, 2023.
- Roma Patel and Ellie Pavlick. Mapping language models to grounded conceptual spaces. In *ICLR*, 2022.
- Michael J. Pazzani and Daniel Billsus. Content-based recommendation systems. In Peter Brusilovsky, Alfred Kobsa, and Wolfgang Nejdl (eds.), *TAW*, 2007.
- Zhaopeng Qiu, Xian Wu, Jingyue Gao, and Wei Fan. U-BERT: pre-training user representations for improved recommendation. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*. AAAI, 2021.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- Shashank Rajput, Nikhil Mehta, Anima Singh, Raghunandan Hulikal Keshavan, Trung Vu, Lukasz Heldt, Lichan Hong, Yi Tay, Vinh Q. Tran, Jonah Samost, Maciej Kula, Ed H. Chi, and Mahesh Sathiamoorthy. Recommender systems with generative retrieval. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *NeurIPS*, 2023.
- Abhilasha Ravichander, Yonatan Belinkov, and Eduard H. Hovy. Probing the probing paradigm: Does probing accuracy entail task relevance? In *EACL*, 2021.
- Xubin Ren, Wei Wei, Lianghao Xia, and Chao Huang. A comprehensive survey on self-supervised learning for recommendation. *arXiv preprint arXiv:2404.03354*, 2024a.
- Xubin Ren, Wei Wei, Lianghao Xia, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao Huang. Representation learning with large language models for recommendation. *CoRR*, abs/2310.15950, 2024b.
- Steffen Rendle. Item recommendation from implicit feedback. In *Recommender Systems Handbook*. Springer US, 2022.
- Steffen Rendle, Christoph Freudenthaler, Zeno Gantner, and Lars Schmidt-Thieme. BPR: bayesian personalized ranking from implicit feedback. *CoRR*, abs/1205.2618, 2012.

- Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. Function vectors in large language models. *CoRR*, abs/2310.15213, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023b.
- Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018.
- Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017.
- Arpita Vats, Vinija Jain, Rahul Raja, and Aman Chadha. Exploring the impact of large language models on recommender systems: An extensive review. *CoRR*, abs/2402.18590, 2024.
- Ivan Vulic, Edoardo Maria Ponti, Robert Litschko, Goran Glavas, and Anna Korhonen. Probing pretrained language models for lexical semantics. In *EMNLP*, 2020.
- Feng Wang and Huaping Liu. Understanding the behaviour of contrastive loss. In *CVPR*, 2021.
- Shoujin Wang, Liang Hu, Yan Wang, Longbing Cao, Quan Z. Sheng, and Mehmet A. Orgun. Sequential recommender systems: Challenges, progress and prospects. In *IJCAI*, 2019a.
- Xiang Wang, Xiangnan He, Meng Wang, Fuli Feng, and Tat-Seng Chua. Neural graph collaborative filtering. In *SIGIR*, 2019b.
- Yuling Wang, Changxin Tian, Binbin Hu, Yanhua Yu, Ziqi Liu, Zhiqiang Zhang, Jun Zhou, Liang Pang, and Xiao Wang. Can small language models be good reasoners for sequential recommendation? *CoRR*, abs/2403.04260, 2024.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, 2022.
- Wei Wei, Xubin Ren, Jiabin Tang, Qinyong Wang, Lixin Su, Suqi Cheng, Junfeng Wang, Dawei Yin, and Chao Huang. Llmrec: Large language models with graph augmentation for recommendation. In Luz Angelica Caudillo-Mata, Silvio Lattanzi, Andrés Muñoz Medina, Leman Akoglu, Aristides Gionis, and Sergei Vassilvitskii (eds.), *WSDM*, 2024.
- Felix Wu, Amauri H. Souza Jr., Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Q. Weinberger. Simplifying graph convolutional networks. In *ICML*, 2019.
- Jiancan Wu, Xiang Wang, Fuli Feng, Xiangnan He, Liang Chen, Jianxun Lian, and Xing Xie. Self-supervised graph learning for recommendation. In *SIGIR*, 2021.

- Jiancan Wu, Xiang Wang, Xingyu Gao, Jiawei Chen, Hongcheng Fu, Tianyu Qiu, and Xiangnan He. On the effectiveness of sampled softmax loss for item recommendation. *CoRR*, abs/2201.02327, 2022.
- Likang Wu, Zhi Zheng, Zhaopeng Qiu, Hao Wang, Hongchao Gu, Tingjia Shen, Chuan Qin, Chen Zhu, Hengshu Zhu, Qi Liu, Hui Xiong, and Enhong Chen. A survey on large language models for recommendation. *CoRR*, abs/2305.19860, 2023a.
- Shiwen Wu, Fei Sun, Wentao Zhang, Xu Xie, and Bin Cui. Graph neural networks in recommender systems: A survey. *ACM Comput. Surv.*, 55(5):97:1–97:37, 2023b.
- Yunjia Xi, Weiwen Liu, Jianghao Lin, Jieming Zhu, Bo Chen, Ruiming Tang, Weinan Zhang, Rui Zhang, and Yong Yu. Towards open-world recommendation with knowledge augmentation from large language models. *CoRR*, abs/2306.10933, 2023.
- Lanling Xu, Junjie Zhang, Bingqian Li, Jinpeng Wang, Mingchen Cai, Wayne Xin Zhao, and Ji-Rong Wen. Prompting large language models for recommender systems: A comprehensive framework and empirical analysis. *CoRR*, abs/2401.04997, 2024.
- Zhengyi Yang, Jiancan Wu, Yanchen Luo, Jizhi Zhang, Yancheng Yuan, An Zhang, Xiang Wang, and Xiangnan He. Large language model can interpret latent space of sequential recommender. *CoRR*, abs/2310.20487, 2023.
- Junliang Yu, Hongzhi Yin, Xin Xia, Tong Chen, Lizhen Cui, and Quoc Viet Hung Nguyen. Are graph augmentations necessary?: Simple graph contrastive learning for recommendation. In *SIGIR*, 2022.
- Junliang Yu, Xin Xia, Tong Chen, Lizhen Cui, Nguyen Quoc Viet Hung, and Hongzhi Yin. Xsimgl: Towards extremely simple graph contrastive learning for recommendation. *IEEE Trans. Knowl. Data Eng.*, 36(2):913–926, 2024.
- Guanghu Yuan, Fajie Yuan, Yudong Li, Beibei Kong, Shujie Li, Lei Chen, Min Yang, Chenyun Yu, Bo Hu, Zang Li, Yu Xu, and Xiaohu Qie. Tenrec: A large-scale multipurpose benchmark dataset for recommender systems. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *NeurIPS*, 2022.
- Zheng Yuan, Fajie Yuan, Yu Song, Youhua Li, Junchen Fu, Fei Yang, Yunzhu Pan, and Yongxin Ni. Where to go next for recommender systems? ID- vs. modality-based recommender models revisited. In *SIGIR*, 2023.
- Jiaqi Zhai, Lucy Liao, Xing Liu, Yueming Wang, Rui Li, Xuan Cao, Leon Gao, Zhaojie Gong, Fangda Gu, Michael He, Yinghai Lu, and Yu Shi. Actions speak louder than words: Trillion-parameter sequential transducers for generative recommendations. *CoRR*, abs/2402.17152, 2024.
- An Zhang, Wenchang Ma, Xiang Wang, and Tat-Seng Chua. Incorporating bias-aware margins into contrastive loss for collaborative filtering. In *NeurIPS*, 2022.
- An Zhang, Leheng Sheng, Zhibo Cai, Xiang Wang, and Tat-Seng Chua. Empowering collaborative filtering with principled adversarial contrastive loss. In *NeurIPS*, 2023a.
- An Zhang, Leheng Sheng, Yuxin Chen, Hao Li, Yang Deng, Xiang Wang, and Tat-Seng Chua. On generative agents in recommendation. *CoRR*, abs/2310.10108, 2023b.
- Chao Zhang, Shiwei Wu, Haoxin Zhang, Tong Xu, Yan Gao, Yao Hu, and Enhong Chen. Notellm: A retrievable large language model for note recommendation. In *WWW*, 2024a.
- Jiaqi Zhang, Yu Cheng, Yongxin Ni, Yunzhu Pan, Zheng Yuan, Junchen Fu, Youhua Li, Jie Wang, and Fajie Yuan. Ninerec: A benchmark dataset suite for evaluating transferable recommendation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024b.
- Junjie Zhang, Yupeng Hou, Ruobing Xie, Wenqi Sun, Julian J. McAuley, Wayne Xin Zhao, Leyu Lin, and Ji-Rong Wen. Agentcf: Collaborative learning with autonomous language agents for recommender systems. *CoRR*, abs/2310.09233, 2023c.

- Junjie Zhang, Ruobing Xie, Yupeng Hou, Wayne Xin Zhao, Leyu Lin, and Ji-Rong Wen. Recommendation as instruction following: A large language model empowered recommendation approach. *CoRR*, abs/2305.07001, 2023d.
- Lingzi Zhang, Xin Zhou, Zhiwei Zeng, and Zhiqi Shen. Are ID embeddings necessary? whitening pre-trained text embeddings for effective sequential recommendation. *CoRR*, abs/2402.10602, 2024c.
- Yang Zhang, Fuli Feng, Jizhi Zhang, Keqin Bao, Qifan Wang, and Xiangnan He. Collm: Integrating collaborative embeddings into large language models for recommendation. *CoRR*, abs/2310.19488, 2023e.
- Hongke Zhao, Songming Zheng, Likang Wu, Bowen Yu, and Jing Wang. LANE: logic alignment of non-tuning large language models and online recommendation systems for explainable reason generation. *CoRR*, abs/2407.02833, 2024.
- Bowen Zheng, Yupeng Hou, Hongyu Lu, Yu Chen, Wayne Xin Zhao, Ming Chen, and Ji-Rong Wen. Adapting large language models by integrating collaborative semantics for recommendation. *CoRR*, abs/2311.09049, 2023.
- Guorui Zhou, Xiaoqiang Zhu, Chengru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. Deep interest network for click-through rate prediction. In *KDD*. ACM, 2018.
- Feng Zhu, Yan Wang, Chaochao Chen, Jun Zhou, Longfei Li, and Guanfeng Liu. Cross-domain recommendation: Challenges, progress, and prospects. In *IJCAI*, 2021.
- Yaochen Zhu, Liang Wu, Qi Guo, Liangjie Hong, and Jundong Li. Collaborative large language model for recommender systems. *CoRR*, abs/2311.01343, 2023.

A RELATED WORKS

Representations in LMs. The impressive capabilities demonstrated by LMs across various tasks raise a wide concern about what they have learned in the representation space. Linear methods (*e.g.*, linear probing (Merullo et al., 2023) and linear mapping (Alain & Bengio, 2017)) are important and effective approaches for interpreting and analyzing representations of LMs (Ravichander et al., 2021). The main idea of linear methods is simple: training linear classifiers to predict some specific attributes or concepts (*e.g.*, lexical structure (Vulic et al., 2020)) from the representations in the hidden layers of LMs, or transforming language representations into another feature space with a linear matrix. A high result of linear methods (*e.g.*, classification accuracy on the out-of-sample test set) tends to imply relevant information has been implicitly encoded in the representation space of LMs, although this does not imply LMs directly use these representations (Ravichander et al., 2021; Gurnee & Tegmark, 2023). Recent studies empirically demonstrate that concepts such as color (Patel & Pavlick, 2022), game states (Li et al., 2023a), and geographic position are encoded in LMs. Furthermore, these concepts may even be linearly encoded in the representation space of LMs (Li et al., 2023a; Park et al., 2023).

Collaborative filtering. Collaborative filtering (CF) (Ren et al., 2024a) is an advanced technique in modern recommender systems. The prevailing CF methods tend to adopt an ID-based paradigm, where users and items are typically represented as one-hot vectors, with an embedding table used for lookup (Koren et al., 2009). Usually, these embedding parameters are learned by optimizing specific loss functions to reconstruct the history interaction pattern (Rendle et al., 2012). Recent advances in CF mainly benefit from two aspects, graph convolution (Wu et al., 2023b) and contrastive learning (Ren et al., 2024a). These CF models exhibit superior recommendation performance by conducting the embedding propagation (Wang et al., 2019b; He et al., 2021) and applying contrastive learning objectives (Wu et al., 2021; Cai et al., 2023; Yu et al., 2024). However, although effective, these methods are still limited, due to the ID-based paradigm. Since one-hot vectors contain no feature information beyond being identifiers, it is challenging to transfer pre-trained ID embeddings to other domains (Hou et al., 2022) or to leverage leading techniques from computer vision (CV) and natural language processing (NLP) (Yuan et al., 2023).

LMs for recommendation. The remarkable language understanding and reasoning ability shown by LMs has attracted extensive attention in the field of recommendation. The application of LMs in recommendation can be categorized into three main approaches: LM-enhanced recommendation, LM as the modality encoder, and LM-as recommender. The first research direction, LM-enhanced recommendation, focuses on empowering traditional recommenders with the semantic representations from LMs (Xi et al., 2023; Ren et al., 2024b; Wei et al., 2024; Geng et al., 2024; Chen, 2023; Wang et al., 2024; Hou et al., 2023; Mao et al., 2023; Qiu et al., 2021; Zhang et al., 2024c). Specifically, these methods introduce representations from LMs as additional features for traditional ID-based recommenders, to capture complicated user preferences. The second research line lies in adopting the LM as the text modality encoder, which is also known as a kind of modality-based recommendation (MoRec) (Yuan et al., 2023; Li et al., 2023c). These methods tend to train the LM as the text modality encoder together with the traditional recommender. In previous studies, BERT-style LMs are widely used as the text modality encoder. The third research line fails in directly using LMs as the recommender and recommends items in a text generation paradigm. Early attempts focus on adopting in-context learning (ICL) (Dong et al., 2022) and prompting pre-trained LMs (Hou et al., 2024b; Liu et al., 2023a; Dai et al., 2023; Gao et al., 2023). However, such naive methods tend to yield poor performance compared to traditional models. Therefore, recent studies concentrate on fine-tuning LMs on recommendation-related corpus (Bao et al., 2023; Zhang et al., 2023d; Lin et al., 2023b; Cui et al., 2022; Liu et al., 2024b; Hua et al., 2023; Chen et al., 2024) and align the LMs with the representations from traditional recommenders as the additional modality (Liao et al., 2024; Zhang et al., 2023e; Yang et al., 2023; Li et al., 2023d; Kong et al., 2024).

B UNCOVERING COLLABORATIVE SIGNALS IN LMS VIA LINEAR MAPPING

B.1 BRIEF OF USED LMS

We briefly introduce the LMs we use for linear mapping in Section 3.1.

Table 6: Dataset statistics.

	Books	Movies & TV	Video Games	Industrial & Scientific	MovieLens-1M	Book Crossing
#Users	71,306	26,073	40,834	15,141	6,040	6,273
#Items	26,073	12,464	14,344	5,163	3,043	5,335
#Interactions	2,209,030	876,027	390,013	82,578	995,492	253,057
Density	0.0008	0.0026	0.0007	0.0010	0.0542	0.0076

- **BERT** (Devlin et al., 2019) is an encoder-only language model based on the transformer architecture (Vaswani et al., 2017), pre-trained on text corpus with unsupervised tasks. BERT adopts bidirectional self-attention heads to learn bidirectional representations.
- **RoBERTa** (Liu et al., 2019) is an enhanced version of BERT. RoBERTa preserves the architecture of BERT but improves it by training with more data and large batches, adopting dynamic masking, and removing the next sentence prediction objective.
- **Llama2-7B** (Touvron et al., 2023b) is an open-source decoder-only LLM with 7 billion parameters. Llama2 adopts grouped-query attention, with longer context length and larger size of the pre-training corpus compared with Llama-7B (Touvron et al., 2023a).
- **Mistral-7B** (Jiang et al., 2023) is an open-source pre-trained decoder-only LLM with 7 billion parameters. Mistral 7B leverages grouped-query attention, coupled with sliding window attention for faster and lower cost inference.
- **text-embedding-ada-v2 & text-embeddings-3-large** (Neelakantan et al., 2022) are leading text embedding models released by OpenAI. These models are built upon decoder-only GPT models, pre-trained on unsupervised data at scale.
- **SFR-Embedding-Mistral** (Meng et al., 2024) is a decoder-based text embedding model built upon the open-source LLM Mixtral-7B (Jiang et al., 2023). SFR-Embedding-Mistral introduces task-homogeneous batching and computes contrastive loss on “hard negatives”, which brings a better performance than the vanilla Mixtral-7B model.

B.2 GENERATING ITEM REPRESENTATIONS FROM LMS

We present how to extract representations from LMs. For encoder-based LMs (e.g., BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019)), we use the representation of the last hidden state corresponding to the [CLS] token (Hou et al., 2024a). For decoder-based models (e.g., Llama-7B (Touvron et al., 2023b; Jiang et al., 2023), Mistral-7B, and SFR-Embedding-Mistral (Meng et al., 2024)), we use the representation in the last transformer block (Vaswani et al., 2017), corresponding to the last input token (Gurnee & Tegmark, 2023; Todd et al., 2023; Neelakantan et al., 2022). Especially, for the commercial closed-source model (e.g., text-embedding-ada-v2 and text-embeddings-3-large² (Neelakantan et al., 2022)), we directly call the API interface to obtain representations.

C LEVERAGING LANGUAGE REPRESENTATIONS FOR BETTER RECOMMENDATION

C.1 DATASETS

We incorporate six datasets in this paper, including four datasets from the Amazon platform³ (Ni et al., 2019) (i.e., Books, Movies & TV, Video Games, and Industrial & Scientific), and two datasets from other platforms (i.e., MovieLens-1M and Book Crossing). Table 6 reports the dataset statistics.

We divide the history interaction of each user into training, validation, and testing sets with a ratio of 4:3:3, and remove users with less than 20 interactions following previous studies (Zhang et al., 2023b). We also remove items from the test and validation sets that do not appear in the training set, to address the cold start problem.

In this paper, we only use the item titles as the text description. Figure 5 gives some item title examples from different datasets.

²<https://platform.openai.com/docs/guides/embeddings>

³www.amazon.com

Item Title Examples
<p>Books: <i>Dismissed with Prejudice: A J.P. Beaumont Novel; Die for Love: A Jacqueline Kirby Novel of Suspense; The Cloud; Memories Before and After the Sound of Music: An Autobiography; Harry Potter and the Sorcerer's Stone;</i></p> <p>Movies & TV: <i>Batman Begins; Fantastic Four; Max Headroom: The Complete Series; Madagascar; Land of the Dead; King Kong;</i></p> <p>Video Games: <i>USB Microphone for RockBand or Guitar Hero (PS3, Wii, Xbox360); Command & Conquer: Tiberian Sun - PC; Tomb Raider III: Adventures of Lara Croft; Kartia: The Word of Fate; Snowboard Kids; Command & Conquer: Tiberian Sun - PC; Final Fantasy VII; Grim Fandango - PC; Half-Life - PC;</i></p> <p>MovieLens-1M: <i>Basquiat (1996); Tin Cup (1996); Godfather, The (1972); Supercop (1992); Manny & Lo (1996); Bound (1996); Carpool (1996);</i></p> <p>Book Crossing: <i>Prague : A Novel; Chocolate Jesus; Wie Barney es sieht; To Kill a Mockingbird; Sturmzeit. Roman; A Soldier of the Great War; Pride and Prejudice (Dover Thrift Editions);</i></p> <p>Industrial & Scientific: <i>Jurassic Perisphinctes Ammonites from France; FS9140: Spinosaurus - Dinosaur Tooth 20-30mm; FS9410: USA Eocene, Fossil Fish (Knightia alt), A-grade; Delta 50-857 Charcoal Filter for 50-868; Hitachi RP30SA 7-1/2 Gallon Stainless Steel Industrial Shop Vacuum (Discontinued by Manufacturer); Makita 632002-4 14-Inch Cut-Off Wheels (5-Pack) (Discontinued by Manufacturer); PORTER-CABLE 740001801 4 1/2-Inch by 10yd 180 Grit Adhesive-Backed Sanding Roll;</i></p>

Figure 5: Example of item titles.

C.2 BASELINES

We incorporate a series of ID-based CF models as our baselines for general recommendation. These models are classified as classical CF methods (MF, MultVAE, and LightGCN), CL-based CF methods (SGL, BC Loss, and XSimGCL), and LM-enhanced CF methods (KAR, RLMRec). For these LM-enhanced CF methods, we adopt the leading method XSimGCL as the backbone.

- **MF** (Koren et al., 2009; Rendle et al., 2012) is the most basic CF model. It denotes users and items with ID-based embeddings and conducts matrix factorization with Bayesian personalized ranking (BPR) loss.
- **MultVAE** (Liang et al., 2018) is a traditional CF model based on the variational autoencoder (VAE). It regards the item recommendation as a generative process from a multinomial distribution and uses variational inference to estimate parameters. We adopt the same model structure as suggested in the paper: $600 \rightarrow 200 \rightarrow 600$.
- **LightGCN** (He et al., 2021) is a light graph convolution network tailored for the recommendation, which deletes redundant feature transformation and activation function in NGCF (Wang et al., 2019b).
- **SGL** (Wu et al., 2021) introduces graph contrastive learning into recommender models for the first time. By employing node or edge dropout to generate augmented graph views and conduct contrastive learning between two views, SGL achieves better performance than LightGCN.
- **BC Loss** (Zhang et al., 2022) introduces a robust and model-agnostic contrastive loss, handling various data biases in recommendation, especially for popularity bias.
- **XSimGCL** (Yu et al., 2024) directly generates augmented views by adding noise into the inner layer of LightGCN without graph augmentation. The simplicity of XSimGCL leads to a faster convergence speed and better performance.
- **KAR** (Xi et al., 2023) enhances recommender models by integrating knowledge from LMs. It generates textual descriptions of users and items and combine the LM representations with traditional recommenders using a hybrid-expert adaptor.
- **RLMRec** (Ren et al., 2024b) aligns semantic representations of users and items with the representations in CF models through a contrastive loss, as an additional loss trained together with the

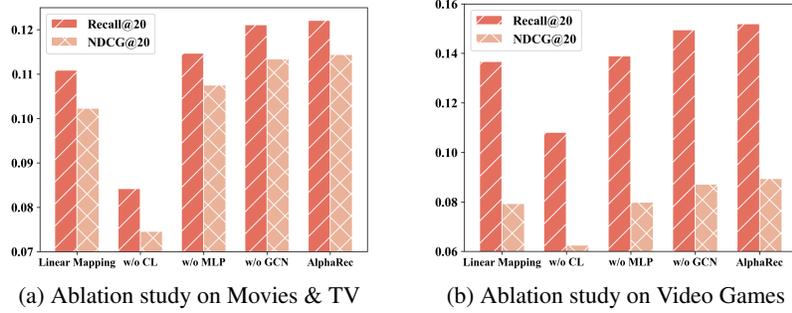


Figure 6: Ablation study

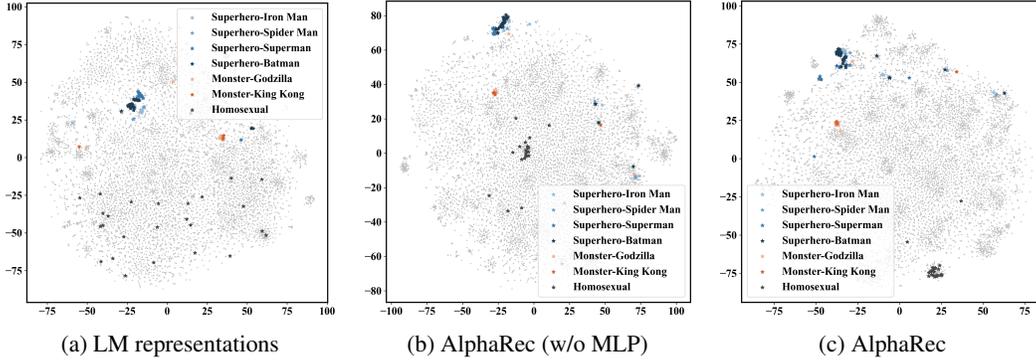


Figure 7: The t-SNE visualization of representations on Movies & TV. (7a) The item representations in the LM space. (7b) The item representations obtained by replacing the MLP with a linear mapping matrix in AlphaRec. (7c) The item representations obtained from AlphaRec.

CF model. The fusion of semantic information and collaborative information brings performance improvement.

C.3 ABLATION STUDY

We conduct the same ablation study as introduced in Section 4 on Movies & TV and Video Games datasets. As illustrated in Figure 6, each component in AlphaRec contributes positively, which is consistent with our findings in Section 4.

C.4 THE T-SNE VISUALIZATION COMPARISON

In this section, we aim to intuitively explore how the MLP in AlphaRec further helps in excavating collaborative signals in language representations, compared to the linear mapping matrix. We visualize the item representations from LMs, post-mapping representations from AlphaRec (w/o MLP), and post-mapping representations from AlphaRec in Figure 7, where AlphaRec (w/o MLP) denotes replacing the MLP with a linear mapping matrix. We observed that movies about superhero and monster cluster in all representation spaces, indicating both AlphaRec (w/o MLP) and AlphaRec capture the preference similarities between these items and preserve the clustering relationship. The difference between AlphaRec (w/o MLP) and AlphaRec lies in the ability to capture obscure preference similarities among items. As shown in Figure 7a, homosexual movies are dispersed in the language space, indicating the possible semantic differences between them. AlphaRec successfully captures the preference similarities and gathers these items in the representation space, while AlphaRec (w/o MLP) remains some items dispersed. Moreover, AlphaRec outperforms AlphaRec (w/o MLP) by a large margin, as indicated in Figure 6a. These results indicate that AlphaRec exhibits a more fine-grained preference capture ability with the help of nonlinear transformation.

Table 7: Training cost of AlphaRec (seconds per epoch/in total).

	Books	Movies & TV	Video Games	Amazon-Mix
AlphaRec	40.1 / 1363.4	12.3 / 479.7	7.4 / 214.6	107.2 / 5788.8

D EXPLORING POTENTIALS OF LANGUAGE REPRESENTATIONS FOR RECOMMENDATION

D.1 FAST CONVERGENCE SPEED

We report the training cost of AlphaRec in this section. Table 7 reports the seconds needed per epoch and the total training cost until convergence. Here Amazon-Mix denotes the mixed dataset of Books, Movies & TV, and Video Games. It’s worth noting that AlphaRec converges quickly and only requires a small amount of training time.

D.2 ZERO-SHOT ABILITY

D.2.1 CO-TRAINING ON MULTIPLE DATASETS

Co-training on multiple datasets is similar to training on one single dataset, where the only difference lies in the negative sampling. When co-training on multiple datasets, the negative items are restricted to the same dataset as the positive item rather than the full item pool. The other training procedures remain the same with training on one single dataset.

D.2.2 BASELINES

Since previous works about zero-shot recommendation mostly focus on sequential recommendation (Kang & McAuley, 2018; Wang et al., 2019a), we slightly modify two methods in sequential recommendation, ZESRec (Ding et al., 2021) and UniSRec (Hou et al., 2022) as our baselines. Specifically, we maintain the model structure as provided in the paper, and adopt the training paradigm of CF.

- **Random** denotes randomly recommending items from the entire item pool.
- **Pop** denotes randomly recommending from the most popular items. Here popularity denotes the number of users that have interacted with the item.
- **ZESRec** (Ding et al., 2021) is the first work that defines the problem of zero-shot recommendation. To address this problem, this work introduces a hierarchical Bayesian model with representations from the pre-trained BERT.
- **UniSRec** (Hou et al., 2022) aims to learn universal item representations from BERT, with parametric whitening and a MoE-enhanced adaptor. By pre-training on multiple source datasets, UniS-Rec can conduct zero-shot recommendation on various datasets in a transductive or inductive paradigm.

D.2.3 COMPARISON WITH LLMRANK

Table 8: Zero-shot performance comparison with LLMRank

	MovieLens-1M				Steam			
	NDCG@1	NDCG@5	NDCG@10	NDCG@20	NDCG@1	NDCG@5	NDCG@10	NDCG@20
LLMRank	0.2485	0.4115	0.5249	0.5612	0.3112	0.4413	0.5255	0.5302
AlphaRec	0.3919	0.6038	0.6543	0.6672	0.4450	0.6131	0.6394	0.6714
Imp. %	57.71%	46.73%	24.65%	18.89%	42.99%	38.93%	21.67%	26.63%

Table 8 illustrates the zero-shot recommendation performance compared with the LLM4Rec method LLMRank. We adopt the same setting of LLMRank, equipping 19 negative items for each positive item, and evaluate the NDCG on the candidate set. AlphaRec exhibits excellent zero-shot performance, significantly surpassing LLMRank. Moreover, the improvement over LLMRank exhibits a rising trend as the K of NDCG decreases.

Table 9: The effect of the training dataset on zero-shot recommendation

	Industrial & Scientific			MovieLens-1M			Book Crossing		
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
AlphaRec (trained on Books)	0.0896	0.0562	0.1256	0.1218	0.2619	0.8942	<u>0.0646</u>	<u>0.0532</u>	<u>0.3346</u>
AlphaRec (trained on Movies & TV)	<u>0.0909</u>	0.0581	<u>0.1266</u>	<u>0.1438</u>	<u>0.3122</u>	<u>0.9200</u>	0.0471	0.0406	0.2600
AlphaRec (trained on Video Games)	0.0905	0.0567	0.1225	0.1221	0.2313	0.9034	0.0412	0.0378	0.2585
AlphaRec (trained on mixed dataset)	0.0913	<u>0.0573</u>	0.1277	0.1486	0.3215	0.9296	0.0660	0.0545	0.3381

Table 10: Performance comparison between training on the single dataset and the mixed dataset

	Books			Movies & TV			Video Games		
	Recall	NDCG	HR	Recall	NDCG	HR	Recall	NDCG	HR
AlphaRec (trained on single dataset)	0.0991	0.0828	0.4185	0.1221	0.1144	0.5587	0.1519	0.0894	0.3207
AlphaRec (trained on mixed dataset)	0.0979	0.0818	0.4147	0.1194	0.1107	0.5463	0.1381	0.0827	0.2985

D.2.4 THE EFFECT OF TRAINING DATASETS

The effect of the training dataset scale on zero-shot recommendation. We report the zero-shot recommendation performance differences trained on different datasets in Table 9. Here AlphaRec (trained on Books) denotes training on a single Books dataset, while AlphaRec (trained on mixed dataset) denotes co-training on three Amazon datasets. Generally, training on more datasets leads to a better zero-shot performance. In addition, we observe that, for the zero-shot performance on untrained target datasets, training datasets with similar themes contribute more (*e.g.*, Movies & TV and MovieLens-1M).

The performance comparison between training on the single dataset and the mixed dataset. In Table 10, AlphaRec (trained on single dataset) denotes training and testing on the same single dataset, while AlphaRec (trained on mixed dataset) denotes training on three Amazon datasets (*i.e.*, Books, Movies & TV, and Video Games) and testing on one single dataset. Generally, co-training on three Amazon datasets yields similar performance compared with training on one single dataset. The only exception lies in Video Games, which shows some performance degradation. We attribute this to the difference between the selection of τ . We use $\tau = 0.15$ when trained on the mixed dataset, while the optimal τ for Video Games lies around 0.2. These results indicate that a single AlphaRec can capture user preferences among various datasets, showcasing a general collaborative signal capture ability.

D.3 INTENTION-AWARE ABILITY

D.3.1 INTENTION QUERY GENERATION

Intention Query Generation

Input
 You are an expert in generating queries for a target movie. Please help me generate the most suitable query for the target movie within one sentence, following the given example.
 Example:
 TARGET: [BUG-A-SALT 3.0 Black Fly Edition](#).
 QUERY: I want a gun that I can use while gardening to get rid of stink bugs, ants, flies, and spiders in my house. It needs to be amazing and help me feel less scared.
 TARGET: [Toy Story \(1995\)](#).

Output
 QUERY: I'm looking for a heartwarming animated movie that follows the adventures of a group of toys who come to life when their owner is not around.

Figure 8: Example of item query generation.

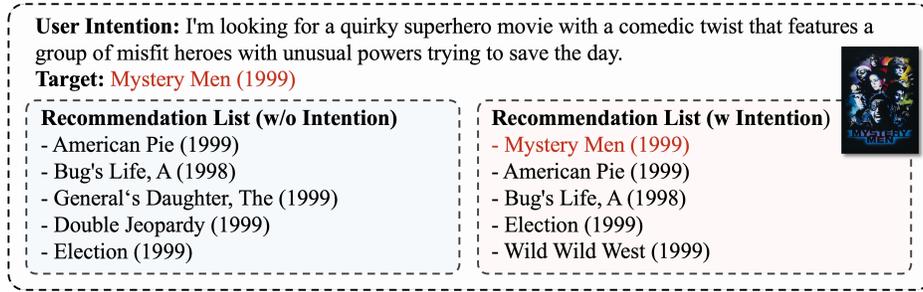


Figure 9: Case study of user intention capture on MovieLens-1M

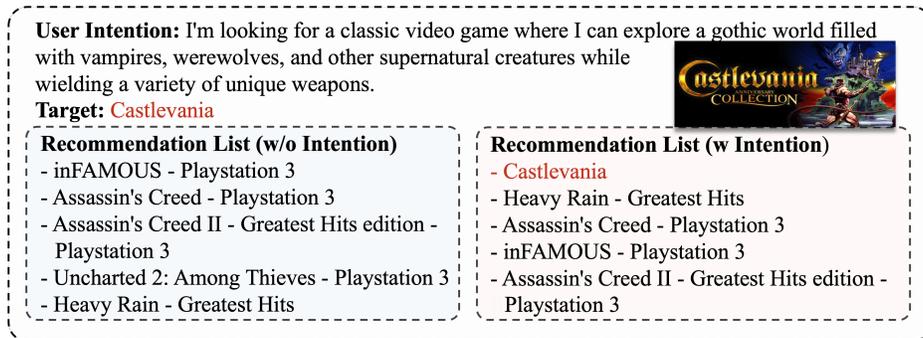


Figure 10: Case study of user intention capture on Video Games

The user intention query is a natural language sentence implying the target item of interest. For each item in the dataset, we generate a fixed user intention query. Following the previous work (Hou et al., 2024a), we generate user intention queries with the help of ChatGPT (OpenAI, 2023). As shown in Figure 8, we prompt ChatGPT in a Chain-of-Thought (CoT) (Wei et al., 2022) paradigm and adopt the output as the user intention query. We adopt a rule-based strategy to ensure the quality of generated queries, and regenerate the wrong query. Considering the huge amount of item title text, we use ChatGPT3.5 API for generating all queries for the budget's sake.

D.3.2 BASELINE

AlphaRec exhibits user intention capture abilities, although not specially designed for search tasks. We compare AlphaRec with TEM (Bi et al., 2020) which falls in the field of personalized search (Ai et al., 2017; McAuley et al., 2015).

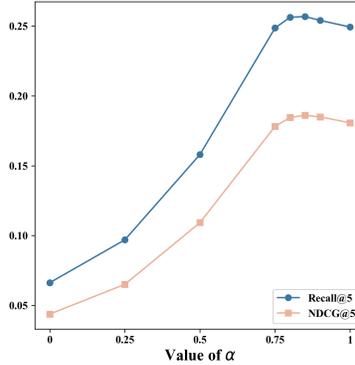
- **TEM** (Bi et al., 2020) uses a transformer to encode the intention query together with user history behaviors, which enables it to achieve better search results by considering the user's historical interest.

D.3.3 CASE STUDY

We conduct two more case studies to verify the user intention capture ability of AlphaRec. As illustrated in Figure 9 and Figure 10, AlphaRec provides better recommendation results, assigning the target item at the top while maintaining the general user preferences.

D.3.4 EFFECT OF THE INTENTION STRENGTH ALPHA

The value of α controls the balance between the user's historical interests and the user intention query. A larger α incorporates more about the user intention while considering less about the user's historical interests. As shown in Figure 11, the effect of α on Video Games shows a similar trend with MovieLens-1M.

Figure 11: Effect of α on Video Games

E HYPERPARAMETER SETTINGS AND IMPLEMENTATION DETAILS

Table 11: Hyperparameters search space for baselines.

	Hyperparameter space
MF & LightGCN	$lr \sim \{1e-5, 3e-5, 5e-5, 1e-4, 3e-4, 5e-4, 1e-3\}$
MultVAE	dropout ratio $\sim \{0, 0.2, 0.5\}$, $\beta \sim \{0.2, 0.4, 0.6, 0.8\}$
SGL	$\tau \sim [0.05, 2]$, $\lambda_1 \sim \{0.005, 0.01, 0.05, 0.1, 0.5, 1.0\}$, $\rho \sim \{0, 0.1, 0.2, 0.3, 0.4, 0.5\}$
BC Loss	$\tau_1 \sim [0.05, 3]$, $\tau_2 \sim [0.05, 2]$
XSimGCL	$\tau \sim [0.05, 2]$, $\epsilon \sim \{0.01, 0.05, 0.1, 0.2, 0.5, 1.0\}$, $\lambda \sim \{0.005, 0.01, 0.05, 0.1, 0.5, 1.0\}$, $l^* = 1$
KAR	No. shared experts $\sim \{3, 4, 5\}$, No. preference experts $\sim \{4, 5\}$
RLMRec	kd weight $\sim [0.05, 2]$, kd temperature $\sim [0.01, 0.05, 0.1, 0.15, 0.2, 0.5, 1]$
ZESRec	$\lambda_u \sim \{0.01, 0.05, 0.1, 0.5, 1.0\}$, $\lambda_v \sim \{0.01, 0.05, 0.1, 0.5, 1.0\}$
UniSRec	$lr \sim \{3e-4, 1e-3, 3e-3, 1e-2\}$
TEM	$l \sim \{2, 3\}$, head $h \sim \{4, 8\}$
AlphaRec	$\tau \sim [0.05, 2]$

Table 12: The hyperparameters of AlphaRec

	Books	Movies & TV	Video Games	Amazon-Mix
τ	0.15	0.15	0.2	0.15

We conduct all the experiments in PyTorch with a single NVIDIA RTX A5000 (24G) GPU and a 64 AMD EPYC 7543 32-Core Processor CPU. We optimize all methods with the Adam optimizer. For all ID-based CF methods, we set the layer numbers of graph propagation by default at 2, with the embedding size as 64 and the size of sampled negative items $|\mathcal{S}_u|$ as 256. We use the early stop strategy to avoid overfitting. We stop the training process if the Recall@20 metric on the validation set does not increase for 20 successive evaluations. In AlphaRec, the dimensions of the input and output in the two-layer MLP are 3072 and 64 respectively, with the hidden layer dimension as 1536. We apply the all-ranking strategy (Krichene & Rendle, 2020) for all experiments, which ranks all items except positive ones in the training set for each user. We search hyperparameters for baselines according to the suggestion in the literature. The hyperparameter search space is reported in Table 11. For these LM-enhanced models, KAR and RLMRec, we also search the hyperparameter of their backbone XSimGCL.

For AlphaRec, the only hyperparameter is the temperature τ and we search it in $[0.05, 2]$. We report the temperature τ we used for each dataset in Table 12. For the mixed dataset Amazon-Mix in Section 5.2, we use a universal $\tau = 0.15$. We adopt $\tau = 0.2$ for the MovieLens-1M dataset for the user intention capture experiment in Section 5.3.